

# Finite difference methods for diffusion processes

Hans Petter Langtangen<sup>1,2</sup>

<sup>1</sup>Center for Biomedical Computing, Simula Research Laboratory

<sup>2</sup>Department of Informatics, University of Oslo

Aug 18, 2014

Note: **VERY PRELIMINARY VERSION**

## Contents

<b>1</b>	<b>The 1D diffusion equation</b>	<b>3</b>
1.1	The initial-boundary value problem for 1D diffusion . . . . .	3
1.2	Forward Euler scheme . . . . .	4
1.3	Backward Euler scheme . . . . .	5
1.4	Sparse matrix implementation . . . . .	8
1.5	The $\theta$ rule . . . . .	8
1.6	The Laplace and Poisson equation . . . . .	9
1.7	Extensions . . . . .	9
<b>2</b>	<b>Analysis of schemes for the diffusion equation</b>	<b>10</b>
2.1	Properties of the solution . . . . .	10
2.2	Analysis of discrete equations . . . . .	11
2.3	Analysis of the finite difference schemes . . . . .	13
2.4	Analysis of the Forward Euler scheme . . . . .	13
2.5	Analysis of the Backward Euler scheme . . . . .	15
2.6	Analysis of the Crank-Nicolson scheme . . . . .	15
2.7	Summary of accuracy of amplification factors . . . . .	16
<b>3</b>	<b>Exercises</b>	<b>19</b>

## List of Exercises and Projects

Exercise	1	Use an analytical solution to formulate a ...	
Exercise	2	Use an analytical solution to formulate a ...	
Exercise	3	Examine stability of a diffusion model with ...	p. 18
Exercise	4	Stabilizing the Crank-Nicolson method by Rannacher ...	p. 19
Project	5	Energy estimates for diffusion problems	p. 19

# 1 The 1D diffusion equation

indexdiffusion equation, 1D indexheat equation, 1D

The famous *diffusion equation*, also known as the *heat equation*, reads

$$\frac{\partial u}{\partial t} = \alpha \frac{\partial^2 u}{\partial x^2},$$

where  $u(x, t)$  is the unknown function to be solved for,  $x$  is a coordinate in space, and  $t$  is time. The coefficient  $\alpha$  is the *diffusion coefficient* and determines how fast  $u$  changes in time. A quick short form for the diffusion equation is  $u_t = \alpha u_{xx}$ .

Compared to the wave equation,  $u_{tt} = c^2 u_{xx}$ , which looks very similar, but the diffusion equation features solutions that are very different from those of the wave equation. Also, the diffusion equation makes quite different demands to the numerical methods.

Typical diffusion problems may experience rapid change in the very beginning, but then the evolution of  $u$  becomes slower and slower. The solution is usually very smooth, and after some time, one cannot recognize the initial shape of  $u$ . This is in sharp contrast to solutions of the wave equation where the initial shape is preserved - the solution is basically a moving initial condition. The standard wave equation  $u_{tt} = c^2 u_{xx}$  has solutions that propagates with speed  $c$  forever, without changing shape, while the diffusion equation converges to a *stationary solution*  $\bar{u}(x)$  as  $t \rightarrow \infty$ . In this limit,  $u_t = 0$ , and  $\bar{u}$  is governed by  $\bar{u}''(x) = 0$ . This stationary limit of the diffusion equation is called the *Laplace equation* and arises in a very wide range of applications throughout the sciences.

It is possible to solve for  $u(x, t)$  using an explicit scheme, but the time step restrictions soon become much less favorable than for an explicit scheme for the wave equation. And of more importance, since the solution  $u$  of the diffusion equation is very smooth and changes slowly, small time steps are not convenient and not required by accuracy as the diffusion process converges to a stationary state.

## 1.1 The initial-boundary value problem for 1D diffusion

To obtain a unique solution of the diffusion equation, or equivalently, to apply numerical methods, we need initial and boundary conditions. The diffusion equation goes with one initial condition  $u(x, 0) = I(x)$ , where  $I$  is a prescribed function. One boundary condition is required at each point on the boundary, which in 1D means that  $u$  must be known,  $u_x$  must be known, or some combination of them.

We shall start with the simplest boundary condition:  $u = 0$ . The complete initial-boundary value diffusion problem in one space dimension can then be specified as

$$\frac{\partial u}{\partial t} = \alpha \frac{\partial^2 u}{\partial x^2}, \quad x \in (0, L), \quad t \in (0, T] \quad (1)$$

$$u(x, 0) = I(x), \quad x \in [0, L] \quad (2)$$

$$u(0, t) = 0, \quad t > 0, \quad (3)$$

$$u(L, t) = 0, \quad t > 0. \quad (4)$$

Equation (1) is known as a one-dimensional *diffusion equation*, also often referred to as a *heat equation*. With only a first-order derivative in time, only one *initial condition* is needed, while the second-order derivative in space leads to a demand for two *boundary conditions*. The parameter  $\alpha$  must be given and is referred to as the *diffusion coefficient*.

Diffusion equations like (1) have a wide range of applications throughout physical, biological, and financial sciences. One of the most common applications is propagation of heat, where  $u(x, t)$  represents the temperature of some substance at point  $x$  and time  $t$ . Section ?? goes into several widely occurring applications.

## 1.2 Forward Euler scheme

The first step in the discretization procedure is to replace the domain  $[0, L] \times [0, T]$  by a set of mesh points. Here we apply equally spaced mesh points

$$x_i = i\Delta x, \quad i = 0, \dots, N_x,$$

and

$$t_n = n\Delta t, \quad n = 0, \dots, N_t.$$

Moreover,  $u_i^n$  denotes the mesh function that approximates  $u(x_i, t_n)$  for  $i = 0, \dots, N_x$  and  $n = 0, \dots, N_t$ . Requiring the PDE (1) to be fulfilled at a mesh point  $(x_i, t_n)$  leads to the equation

$$\frac{\partial}{\partial t} u(x_i, t_n) = \alpha \frac{\partial^2}{\partial x^2} u(x_i, t_n), \quad (5)$$

The next step is to replace the derivatives by finite difference approximations. The computationally simplest method arises from using a forward difference in time and a central difference in space:

$$[D_t^+ u = \alpha D_x D_x u]_i^n. \quad (6)$$

Written out,

$$\frac{u_i^{n+1} - u_i^n}{\Delta t} = \alpha \frac{u_{i+1}^n - 2u_i^n + u_{i-1}^n}{\Delta x^2}. \quad (7)$$

We have turned the PDE into algebraic equations, also often called discrete equations. The key property of the equations is that they are algebraic, which makes them easy to solve. As usual, we anticipate that  $u_i^n$  is already computed such that  $u_i^{n+1}$  is the only unknown in (7). Solving with respect to this unknown is easy:

$$u_i^{n+1} = u_i^n + \alpha \frac{\Delta t}{\Delta x^2} (u_{i+1}^n - 2u_i^n + u_{i-1}^n). \quad (8)$$

The computational algorithm then becomes

1. compute  $u_i^0 = I(x_i)$  for  $i = 0, \dots, N_x$
2. for  $n = 0, 1, \dots, N_t$ :
  - (a) apply (8) for all the internal spatial points  $i = 1, \dots, N_x - 1$
  - (b) set the boundary values  $u_i^{n+1} = 0$  for  $i = 0$  and  $i = N_x$

The algorithm is compactly fully specified in Python:

```

x = linspace(0, L, Nx+1) # mesh points in space
dx = x[1] - x[0]
t = linspace(0, T, Nt+1) # mesh points in time
dt = t[1] - t[0]
Fo = a*dt/dx**2
u = zeros(Nx+1)
u_1 = zeros(Nx+1)

# Set initial condition u(x,0) = I(x)
for i in range(0, Nx+1):
    u_1[i] = I(x[i])

for n in range(0, Nt):
    # Compute u at inner mesh points
    for i in range(1, Nx):
        u[i] = u_1[i] + Fo*(u_1[i-1] - 2*u_1[i] + u_1[i+1])

    # Insert boundary conditions
    u[0] = 0; u[Nx] = 0

    # Update u_1 before next step
    u_1[:] = u

```

### 1.3 Backward Euler scheme

We now apply a backward difference in time in (5), but the same central difference in space:

$$[D_t^- u = D_x D_x u]_i^n, \quad (9)$$

which written out reads

$$\frac{u_i^n - u_i^{n-1}}{\Delta t} = \alpha \frac{u_{i+1}^n - 2u_i^n + u_{i-1}^n}{\Delta x^2}. \quad (10)$$

Now we assume  $u_i^{n-1}$  is computed, but all quantities at the "new" time level  $n$  are unknown. This time it is not possible to solve with respect to  $u_i^n$  because this value couples to its neighbors in space,  $u_{i-1}^n$  and  $u_{i+1}^n$ , which are also unknown. Let us examine this fact for the case when  $N_x = 3$ . Equation (10) written for  $i = 1, \dots, N_x - 1 = 1, 2$  becomes

$$\frac{u_1^n - u_1^{n-1}}{\Delta t} = \alpha \frac{u_2^n - 2u_1^n + u_0^n}{\Delta x^2} \quad (11)$$

$$\frac{u_2^n - u_2^{n-1}}{\Delta t} = \alpha \frac{u_3^n - 2u_2^n + u_1^n}{\Delta x^2} \quad (12)$$

The boundary values  $u_0^n$  and  $u_3^n$  are known as zero. Collecting the unknown new values  $u_1^n$  and  $u_2^n$  on the left-hand side gives

$$\left(1 + 2\alpha \frac{\Delta t}{\Delta x^2}\right) u_1^n - \alpha \frac{\Delta t}{\Delta x^2} u_2^n = u_1^{n-1}, \quad (13)$$

$$-\alpha \frac{\Delta t}{\Delta x^2} u_1^n + \left(1 + 2\alpha \frac{\Delta t}{\Delta x^2}\right) u_2^n = u_2^{n-1}. \quad (14)$$

This is a coupled  $2 \times 2$  system of algebraic equations for the unknowns  $u_1^n$  and  $u_2^n$ . Discretization methods that lead to a coupled system of equations for the unknown function at a new time

level are said to be *implicit methods*. The counterpart, *explicit methods*, refers to discretization methods where there is a simple explicit formula for the values of the unknown function at each of the spatial mesh points at the new time level. From an implementational point of view, implicit methods are more comprehensive to code since they require the solution of coupled equations, i.e., a matrix system, at each time level.

In the general case, (10) gives rise to a coupled  $(N_x - 1) \times (N_x - 1)$  system of algebraic equations for all the unknown  $u_i^n$  at the interior spatial points  $i = 1, \dots, N_x - 1$ . Collecting the unknowns on the left-hand side, and introducing the *numerical Fourier number*

$$F_o = \alpha \frac{\Delta t}{\Delta x^2}, \quad (15)$$

(10) can be written

$$-F_o u_{i-1}^n + (1 + 2F_o) u_i^n - F_o u_{i+1}^n = u_{i-1}^{n-1}, \quad (16)$$

for  $i = 1, \dots, N_x - 1$ . One can either view these equations as a system for where the  $u_i^n$  values at the internal grid points,  $i = 1, \dots, N_x - 1$ , are unknown, or we may append the boundary values  $u_0^n$  and  $u_{N_x}^n$  to the system. In the latter case, all  $u_i^n$  for  $i = 0, \dots, N_x$  are unknown and we must add the boundary equations to the  $N_x - 1$  equations in (16):

$$u_0^n = 0, \quad (17)$$

$$u_{N_x}^n = 0. \quad (18)$$

A coupled system of algebraic equations can be written on matrix form, and this is important if we want to call up ready-made software for solving the system. The equations (16) and (17)–(18) correspond to the matrix equation

$$AU = b$$

where  $U = (u_0^n, \dots, u_{N_x}^n)$ , and the matrix  $A$  has the following structure:

$$A = \begin{pmatrix} A_{0,0} & A_{0,1} & 0 & \cdots & \cdots & \cdots & \cdots & \cdots & 0 \\ A_{1,0} & A_{1,1} & 0 & \ddots & & & & & \vdots \\ 0 & A_{2,1} & A_{2,2} & A_{2,3} & \ddots & & & & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & 0 & & & \vdots \\ \vdots & & \ddots & \ddots & \ddots & \ddots & \ddots & & \vdots \\ \vdots & & & 0 & A_{i,i-1} & A_{i,i} & A_{i,i+1} & \ddots & \vdots \\ \vdots & & & & \ddots & \ddots & \ddots & \ddots & 0 \\ \vdots & & & & & \ddots & \ddots & \ddots & A_{N_x-1,N_x} \\ 0 & \cdots & \cdots & \cdots & \cdots & \cdots & 0 & A_{N_x,N_x-1} & A_{N_x,N_x} \end{pmatrix} \quad (19)$$

The nonzero elements are given by

$$A_{i,i-1} = -F_o \quad (20)$$

$$A_{i,i} = 1 + 2F_o \quad (21)$$

$$A_{i,i+1} = -F_o \quad (22)$$

for the equations for internal points,  $i = 1, \dots, N_x - 1$ . The equations for the boundary points correspond to

$$A_{0,0} = 1, \quad (23)$$

$$A_{0,1} = 0, \quad (24)$$

$$A_{N_x, N_x - 1} = 0, \quad (25)$$

$$A_{N_x, N_x} = 1. \quad (26)$$

The right-hand side  $b$  is written as

$$b = \begin{pmatrix} b_0 \\ b_1 \\ \vdots \\ b_i \\ \vdots \\ b_{N_x} \end{pmatrix} \quad (27)$$

with

$$b_0 = 0, \quad (28)$$

$$b_i = u_i^{n-1}, \quad i = 1, \dots, N_x - 1, \quad (29)$$

$$b_{N_x} = 0. \quad (30)$$

We observe that the matrix  $A$  contains quantities that do not change in time. Therefore,  $A$  can be formed once and for all before we enter the recursive formulas for the time evolution. The right-hand side  $b$ , however, must be updated at each time step. This leads to the following computational algorithm, here sketched with Python code:

```
x = linspace(0, L, Nx+1) # mesh points in space
dx = x[1] - x[0]
t = linspace(0, T, Nt+1) # mesh points in time
u = zeros(Nx+1)
u_1 = zeros(Nx+1)

# Data structures for the linear system
A = zeros((Nx+1, Nx+1))
b = zeros(Nx+1)

for i in range(1, Nx):
    A[i,i-1] = -Fo
    A[i,i+1] = -Fo
    A[i,i] = 1 + 2*Fo
A[0,0] = A[Nx,Nx] = 1

# Set initial condition u(x,0) = I(x)
for i in range(0, Nx+1):
    u_1[i] = I(x[i])

import scipy.linalg

for n in range(0, Nt):
    # Compute b and solve linear system
```

```
for i in range(1, Nx):
    b[i] = -u_1[i]
b[0] = b[Nx] = 0
u[:] = scipy.linalg.solve(A, b)

# Update u_1 before next step
u_1[:] = u
```

## 1.4 Sparse matrix implementation

We have seen from (19) that the matrix  $A$  is tridiagonal. The code segment above used a full, dense matrix representation of  $A$ , which stores a lot of values we know are zero beforehand, and worse, the solution algorithm computes with all these zeros. With  $N_x + 1$  unknowns, the work by the solution algorithm is  $\frac{1}{3}(N_x + 1)^3$  and the storage requirements  $(N_x + 1)^2$ . By utilizing the fact that  $A$  is tridiagonal and employing corresponding software tools, the work and storage demands can be proportional to  $N_x$  only.

The key idea is to apply a data structure for a tridiagonal or sparse matrix. The `scipy.sparse` package has relevant utilities. For example, we can store the nonzero diagonals of a matrix. The package also has linear system solvers that operate on sparse matrix data structures. The code below illustrates how we can store only the main diagonal and the upper and lower diagonals.

```
# Representation of sparse matrix and right-hand side
main = zeros(Nx+1)
lower = zeros(Nx-1)
upper = zeros(Nx-1)
b = zeros(Nx+1)

# Precompute sparse matrix
main[:] = 1 + 2*Fo
lower[:] = -Fo #1
upper[:] = -Fo #1
# Insert boundary conditions
main[0] = 1
main[Nx] = 1

A = scipy.sparse.diags(
    diagonals=[main, lower, upper],
    offsets=[0, -1, 1], shape=(Nx+1, Nx+1),
    format='csr')
print A.todense()

# Set initial condition
for i in range(0, Nx+1):
    u_1[i] = I(x[i])

for n in range(0, Nt):
    b = u_1
    b[0] = b[-1] = 0.0 # boundary conditions
    u[:] = scipy.sparse.linalg.spsolve(A, b)
    u_1[:] = u
```

The `scipy.sparse.linalg.spsolve` function utilizes the sparse storage structure of  $A$  and performs in this case a very efficient Gaussian elimination solve.

## 1.5 The $\theta$ rule

The  $\theta$  rule provides a family of finite difference approximations in time:

- $\theta = 0$  gives the Forward Euler scheme in time
- $\theta = 1$  gives the Backward Euler scheme in time
- $\theta = \frac{1}{2}$  gives the Crank-Nicolson scheme in time

Applied to the 1D diffusion problem we have

$$\frac{u_i^{n+1} - u_i^n}{\Delta t} = \alpha \left( \theta \frac{u_{i+1}^{n+1} - 2u_i^{n+1} + u_{i-1}^{n+1}}{\Delta x^2} + (1 - \theta) \frac{u_{i+1}^n - 2u_i^n + u_{i-1}^n}{\Delta x^2} \right).$$

This scheme also leads to a matrix system with entries  $1 + 2F_o\theta$  on the main diagonal of the matrix, and  $-F_o\theta$  on the super- and sub-diagonal. The right-hand side entry  $b_i$  is

$$b_i = u_i^n + F_o(1 - \theta) \frac{u_{i+1}^n - 2u_i^n + u_{i-1}^n}{\Delta x^2}.$$

## 1.6 The Laplace and Poisson equation

The Laplace equation,  $\nabla^2 u = 0$ , or the Poisson equation,  $-\nabla^2 u = f$ , occur in numerous applications throughout science and engineering. We can solve 1D variants of the Laplace equations with the listed software, because we can interpret  $u_{xx} = 0$  as the limiting solution of  $u_t = \alpha u_{xx}$  when  $u$  reach a steady state limit where  $u_t \rightarrow 0$ . Similarly, Poisson's equation  $-u_{xx} = f$  arises from solving  $u_t = u_{xx} + f$  and letting  $t \rightarrow \infty$  so  $u_t \rightarrow 0$ .

Technically in a program, we can simulate  $t \rightarrow \infty$  by just taking one large time step, or equivalently, set  $\alpha$  to a large value. All we need is to have  $F_o$  large. As  $F_o \rightarrow \infty$ , we can from the schemes see that the limiting discrete equation becomes

$$\frac{u_{i+1}^{n+1} - 2u_i^{n+1} + u_{i-1}^{n+1}}{\Delta x^2} = 0,$$

which is nothing but the discretization  $[D_x D_x u]_i^{n+1} = 0$  of  $u_{xx} = 0$ .

The Backward Euler scheme can solve the limit equation directly and hence produce a solution of the 1D Laplace equation. With the Forward Euler scheme we must do the time stepping since  $F_o > 1/2$  is illegal and leads to instability. We may interpret this time stepping as solving the equation system from  $u_{xx}$  by iterating on a time pseudo time variable.

## 1.7 Extensions

These extensions are performed exactly as for a wave equation as they only affect the spatial derivatives (which are the same as in the wave equation).

- Variable coefficients
- Neumann and Robin conditions
- 2D and 3D

Future versions of this document will for completeness and independence of the wave equation document feature info on the three points. The Robin condition is new, but straightforward to handle:

$$-\alpha \frac{\partial u}{\partial n} = h_T(u - U_s), \quad [-\alpha D_x u = h_T(u - U_s)]_i^n$$

## 2 Analysis of schemes for the diffusion equation

### 2.1 Properties of the solution

A particular characteristic of diffusive processes, governed by an equation like

$$u_t = \alpha u_{xx}, \tag{31}$$

is that the initial shape  $u(x, 0) = I(x)$  spreads out in space with time, along with a decaying amplitude. Three different examples will illustrate the spreading of  $u$  in space and the decay in time.

**Similarity solution.** The diffusion equation (31) admits solutions that depend on  $\eta = (x - c)/\sqrt{4\alpha t}$  for a given value of  $c$ . One particular solution is

$$u(x, t) = a \operatorname{erf}(\eta) + b, \tag{32}$$

where

$$\operatorname{erf}(\eta) = \frac{2}{\sqrt{\pi}} \int_0^\eta e^{-\zeta^2} d\zeta, \tag{33}$$

is the *error function*, and  $a$  and  $b$  are arbitrary constants. The error function lies in  $(-1, 1)$ , is odd around  $\eta = 0$ , and goes relatively quickly to  $\pm 1$ :

$$\begin{aligned} \lim_{\eta \rightarrow -\infty} \operatorname{erf}(\eta) &= -1, \\ \lim_{\eta \rightarrow \infty} \operatorname{erf}(\eta) &= 1, \\ \operatorname{erf}(\eta) &= -\operatorname{erf}(-\eta), \\ \operatorname{erf}(0) &= 0, \\ \operatorname{erf}(2) &= 0.99532227, \\ \operatorname{erf}(3) &= 0.99997791. \end{aligned}$$

As  $t \rightarrow 0$ , the error function approaches a step function centered at  $x = c$ . For a diffusion problem posed on the unit interval  $[0, 1]$ , we may choose the step at  $x = 1/2$  (meaning  $c = 1/2$ ),  $a = -1/2$ ,  $b = 1/2$ . Then

$$u(x, t) = \frac{1}{2} \left( 1 - \operatorname{erf} \left( \frac{x - \frac{1}{2}}{\sqrt{4\alpha t}} \right) \right) = \frac{1}{2} \operatorname{erfc} \left( \frac{x - \frac{1}{2}}{\sqrt{4\alpha t}} \right), \tag{34}$$

where we have introduced the *complementary error function*  $\operatorname{erfc}(\eta) = 1 - \operatorname{erf}(\eta)$ . The solution (34) implies the boundary conditions

$$u(0, t) = \frac{1}{2} \left( 1 - \operatorname{erf} \left( \frac{-1/2}{\sqrt{4\alpha t}} \right) \right), \tag{35}$$

$$u(1, t) = \frac{1}{2} \left( 1 - \operatorname{erf} \left( \frac{1/2}{\sqrt{4\alpha t}} \right) \right). \tag{36}$$

For small enough  $t$ ,  $u(0, t) \approx 1$  and  $u(1, t) \approx 1$ , but as  $t \rightarrow \infty$ ,  $u(x, t) \rightarrow 1/2$  on  $[0, 1]$ .

**Solution for a Gaussian pulse.** The standard diffusion equation  $u_t = \alpha u_{xx}$  admits a Gaussian function as solution:

$$u(x, t) = \frac{1}{\sqrt{4\pi\alpha t}} \exp\left(-\frac{(x-c)^2}{4\alpha t}\right). \quad (37)$$

At  $t = 0$  this is a Dirac delta function, so for computational purposes one must start to view the solution at some time  $t = t_\epsilon > 0$ . Replacing  $t$  by  $t_\epsilon + t$  in (37) makes it easy to operate with a (new)  $t$  that starts at  $t = 0$  with an initial condition with a finite width. The important feature of (37) is that the standard deviation  $\sigma$  of a sharp initial Gaussian pulse increases in time according to  $\sigma = \sqrt{2\alpha t}$ , making the pulse diffuse and flatten out.

**Solution for a sine component.** For example, (31) admits a solution of the form

$$u(x, t) = Qe^{-at} \sin(kx). \quad (38)$$

The parameters  $Q$  and  $k$  can be freely chosen, while inserting (38) in (31) gives the constraint

$$a = -\alpha k^2.$$

A very important feature is that the initial shape  $I(x) = Q \sin kx$  undergoes a damping  $\exp(-\alpha k^2 t)$ , meaning that rapid oscillations in space, corresponding to large  $k$ , are very much faster dampened than slow oscillations in space, corresponding to small  $k$ . This feature leads to a smoothing of the initial condition with time.

The following examples illustrates the damping properties of (38). We consider the specific problem

$$\begin{aligned} u_t &= u_{xx}, \quad x \in (0, 1), \quad t \in (0, T], \\ u(0, t) &= u(1, t) = 0, \quad t \in (0, T], \\ u(x, 0) &= \sin(\pi x) + 0.1 \sin(100\pi x). \end{aligned}$$

The initial condition has been chosen such that adding two solutions like (38) constructs an analytical solution to the problem:

$$u(x, t) = e^{-\pi^2 t} \sin(\pi x) + 0.1 e^{-\pi^2 10^4 t} \sin(100\pi x). \quad (39)$$

Figure 1 illustrates the rapid damping of rapid oscillations  $\sin(100\pi x)$  and the very much slower damping of the slowly varying  $\sin(\pi x)$  term. After about  $t = 0.5 \cdot 10^{-4}$  the rapid oscillations do not have a visible amplitude, while we have to wait until  $t \sim 0.5$  before the amplitude of the long wave  $\sin(\pi x)$  becomes very small.

## 2.2 Analysis of discrete equations

A counterpart to (38) is the complex representation of the same function:

$$u(x, t) = Qe^{-at} e^{ikx},$$

where  $i = \sqrt{-1}$  is the imaginary unit. We can add such functions, often referred to as wave components, to make a Fourier representation of a general solution of the diffusion equation:

$$u(x, t) \approx \sum_{k \in K} b_k e^{-\alpha k^2 t} e^{ikx}, \quad (40)$$

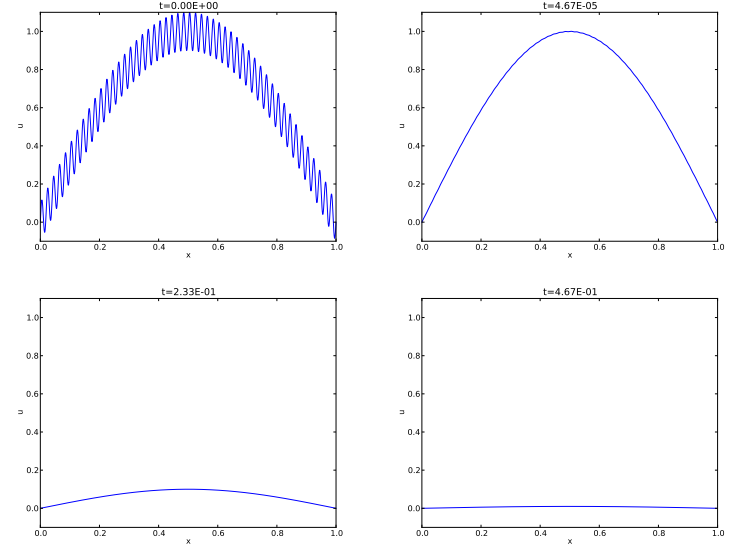


Figure 1: Evolution of the solution of a diffusion problem: initial condition (upper left), 1/100 reduction of the small waves (upper right), 1/10 reduction of the long wave (lower left), and 1/100 reduction of the long wave (lower right).

where  $K$  is a set of an infinite number of  $k$  values needed to construct the solution. In practice, however, the series is truncated and  $K$  is a finite set of  $k$  values need build a good approximate solution. Note that (39) is a special case of (40) where  $K = \{\pi, 100\pi\}$ ,  $b_\pi = 1$ , and  $b_{100\pi} = 0.1$ .

The amplitudes  $b_k$  of the individual Fourier waves must be determined from the initial condition. At  $t = 0$  we have  $u \approx \sum_k b_k \exp(ikx)$  and find  $K$  and  $b_k$  such that

$$I(x) \approx \sum_{k \in K} b_k e^{ikx}. \quad (41)$$

(The relevant formulas for  $b_k$  come from Fourier analysis, or equivalently, a least-squares method for approximating  $I(x)$  in a function space with basis  $\exp(ikx)$ .)

Much insight about the behavior of numerical methods can be obtained by investigating how a wave component  $\exp(-\alpha k^2 t) \exp(ikx)$  is treated by the numerical scheme. It appears that such wave components are also solutions of the schemes, but the damping factor  $\exp(-\alpha k^2 t)$  varies among the schemes. To ease the forthcoming algebra, we write the damping factor as  $A^n$ . The exact amplification factor corresponding to  $A$  is  $A_e = \exp(-\alpha k^2 \Delta t)$ .

### 2.3 Analysis of the finite difference schemes

We have seen that a general solution of the diffusion equation can be built as a linear combination of basic components

$$e^{-\alpha k^2 t} e^{ikx}.$$

A fundamental question is whether such components are also solutions of the finite difference schemes. This is indeed the case, but the amplitude  $\exp(-\alpha k^2 t)$  might be modified (which also happens when solving the ODE counterpart  $u' = -\alpha u$ ). We therefore look for numerical solutions of the form

$$u_q^n = A^n e^{ikq\Delta x} = A^n e^{ikx}, \quad (42)$$

where the amplification factor  $A$  must be determined by inserting the component into an actual scheme.

**Stability.** The exact amplification factor is  $A_e = \exp(-\alpha^2 k^2 \Delta t)$ . We should therefore require  $|A| < 1$  to have a decaying numerical solution as well. If  $-1 \leq A < 0$ ,  $A^n$  will change sign from time level to time level, and we get stable, non-physical oscillations in the numerical solutions that are not present in the exact solution.

**Accuracy.** To determine how accurately a finite difference scheme treats one wave component (42), we see that the basic deviation from the exact solution is reflected in how well  $A^n$  approximates  $A_e^n$ , or how well  $A$  approximates  $A_e$ .

### 2.4 Analysis of the Forward Euler scheme

The Forward Euler finite difference scheme for  $u_t = \alpha u_{xx}$  can be written as

$$[D_t^+ u = \alpha D_x D_x u]_q^n.$$

Inserting a wave component (42) in the scheme demands calculating the terms

$$e^{ikq\Delta x} [D_t^+ A]^n = e^{ikq\Delta x} A^n \frac{A-1}{\Delta t},$$

and

$$A^n D_x D_x [e^{ikx}]_q = A^n \left( -e^{ikq\Delta x} \frac{4}{\Delta x^2} \sin^2 \left( \frac{k\Delta x}{2} \right) \right).$$

Inserting these terms in the discrete equation and dividing by  $A^n e^{ikq\Delta x}$  leads to

$$\frac{A-1}{\Delta t} = -\alpha \frac{4}{\Delta x^2} \sin^2 \left( \frac{k\Delta x}{2} \right),$$

and consequently

$$A = 1 - 4F_o \sin^2 \left( \frac{k\Delta x}{2} \right), \quad (43)$$

where

$$F_o = \frac{\alpha \Delta t}{\Delta x^2} \quad (44)$$

is the *numerical Fourier number*. The complete numerical solution is then

$$u_q^n = \left( 1 - 4F_o \sin^2 \left( \frac{k\Delta x}{2} \right) \right)^n e^{ikq\Delta x}. \quad (45)$$

**Stability.** We easily see that  $A \leq 1$ . However, the  $A$  can be less than  $-1$ , which will lead to growth of a numerical wave component. The criterion  $A \geq -1$  implies

$$4F_o \sin^2(p/2) \leq 2.$$

The worst case is when  $\sin^2(p/2) = 1$ , so a sufficient criterion for stability is

$$F_o \leq \frac{1}{2}, \quad (46)$$

or expressed as a condition on  $\Delta t$ :

$$\Delta t \leq \frac{\Delta x^2}{2\alpha}. \quad (47)$$

Note that halving the spatial mesh size,  $\Delta x \rightarrow \frac{1}{2}\Delta x$ , requires  $\Delta t$  to be reduced by a factor of 1/4. The method hence becomes very expensive for fine spatial meshes.

**Accuracy.** Since  $A$  is expressed in terms of  $F_o$  and the parameter we now call  $p = k\Delta x/2$ , we also express  $A_e$  by  $F_o$  and  $p$ :

$$A_e = \exp(-\alpha k^2 \Delta t) = \exp(-4F_o p^2).$$

Computing the Taylor series expansion of  $A/A_e$  in terms of  $F_o$  can easily be done with aid of *sympy*:

```
def A_exact(Fo, p):
    return exp(-4*Fo*p**2)

def A_FE(Fo, p):
    return 1 - 4*Fo*sin(p)**2

from sympy import *
Fo, p = symbols('Fo p')
A_err_FE = A_FE(Fo, p)/A_exact(Fo, p)
print A_err_FE.series(Fo, 0, 6)
```

The result is

$$\frac{A}{A_e} = 1 - 4F_o \sin^2 p + 2F_o^2 p^2 - 16F_o^2 p^2 \sin^2 p + 8F_o^2 p^4 + \dots$$

Recalling that  $F_o = \alpha \Delta t / \Delta x$ ,  $p = k\Delta x/2$ , and that  $\sin^2 p \leq 1$ , we realize that the dominating error terms are at most

$$1 - 4\alpha \frac{\Delta t}{\Delta x^2} + \alpha \Delta t - 4\alpha^2 \Delta t^2 + \alpha^2 \Delta t^2 \Delta x^2 + \dots$$

## 2.5 Analysis of the Backward Euler scheme

Discretizing  $u_t = \alpha u_{xx}$  by a Backward Euler scheme,

$$[D_t^- u = \alpha D_x D_x u]_q^n,$$

and inserting a wave component (42), leads to calculations similar to those arising from the Forward Euler scheme, but since

$$e^{ikq\Delta x} [D_t^- A]^n = A^n e^{ikq\Delta x} \frac{1 - A^{-1}}{\Delta t},$$

we get

$$\frac{1 - A^{-1}}{\Delta t} = -\alpha \frac{4}{\Delta x^2} \sin^2 \left( \frac{k\Delta x}{2} \right),$$

and then

$$A = (1 + 4F_o \sin^2 p)^{-1}. \quad (48)$$

The complete numerical solution can be written

$$u_q^n = (1 + 4F_o \sin^2 p)^{-n} e^{ikq\Delta x}. \quad (49)$$

**Stability.** We see from (48) that  $0 < A < 1$ , which means that all numerical wave components are stable and non-oscillatory for any  $\Delta t > 0$ .

## 2.6 Analysis of the Crank-Nicolson scheme

The Crank-Nicolson scheme can be written as

$$[D_t u = \alpha D_x D_x \bar{u}]_q^{n+\frac{1}{2}},$$

or

$$[D_t u]_q^{n+\frac{1}{2}} = \frac{1}{2} \alpha ([D_x D_x u]_q^n + [D_x D_x u]_q^{n+1}).$$

Inserting (42) in the time derivative approximation leads to

$$[D_t A^n e^{ikq\Delta x}]^{n+\frac{1}{2}} = A^{n+\frac{1}{2}} e^{ikq\Delta x} \frac{A^{\frac{1}{2}} - A^{-\frac{1}{2}}}{\Delta t} = A^n e^{ikq\Delta x} \frac{A - 1}{\Delta t}.$$

Inserting (42) in the other terms and dividing by  $A^n e^{ikq\Delta x}$  gives the relation

$$\frac{A - 1}{\Delta t} = -\frac{1}{2} \alpha \frac{4}{\Delta x^2} \sin^2 \left( \frac{k\Delta x}{2} \right) (1 + A),$$

and after some more algebra,

$$A = \frac{1 - 2F_o \sin^2 p}{1 + 2F_o \sin^2 p}. \quad (50)$$

The exact numerical solution is hence

$$u_q^n = \left( \frac{1 - 2F_o \sin^2 p}{1 + 2F_o \sin^2 p} \right)^n e^{ikp\Delta x}. \quad (51)$$

**Stability.** The criteria  $A > -1$  and  $A < 1$  are fulfilled for any  $\Delta t > 0$ .

## 2.7 Summary of accuracy of amplification factors

We can plot the various amplification factors against  $p = k\Delta x/2$  for different choices of the  $F_o$  parameter. Figures 2, 3, and 4 show how long and small waves are damped by the various schemes compared to the exact damping. As long as all schemes are stable, the amplification factor is positive, except for Crank-Nicolson when  $F_o > 0.5$ .

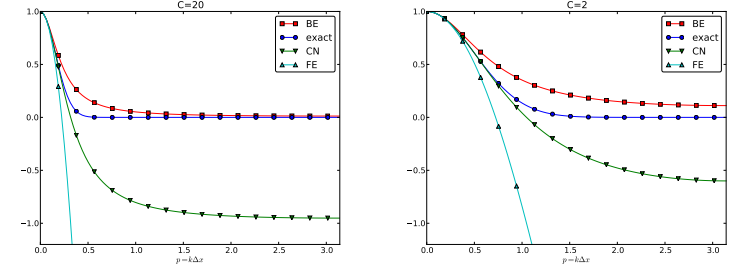


Figure 2: Amplification factors for large time steps.

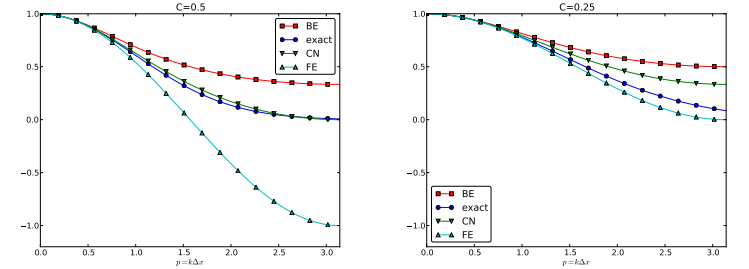


Figure 3: Amplification factors for time steps around the Forward Euler stability limit.

The effect of negative amplification factors is that  $A^n$  changes sign from one time level to the next, thereby giving rise to oscillations in time in an animation of the solution. We see from Figure 2 that for  $F_o = 20$ , waves with  $p \geq \pi/2$  undergo a damping close to  $-1$ , which means that the amplitude does not decay and that the wave component jumps up and down in time. For  $F_o = 2$  we have a damping of a factor of 0.5 from one time level to the next, which is very much smaller than the exact damping. Short waves will therefore fail to be effectively dampened. These waves will manifest themselves as high frequency oscillatory noise in the solution.

A value  $p = \pi/4$  corresponds to four mesh points per wave length of  $e^{ikx}$ , while  $p = \pi/2$  implies only two points per wave length, which is the smallest number of points we can have to represent the wave on the mesh.



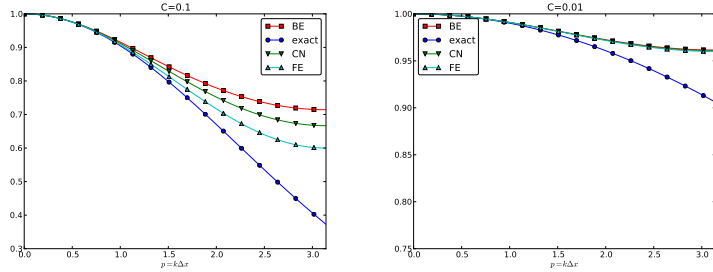


Figure 4: Amplification factors for small time steps.

To demonstrate the oscillatory behavior of the Crank-Nicolson scheme, we choose an initial condition that leads to short waves with significant amplitude. A discontinuous  $I(x)$  will in particular serve this purpose.

Run  $F_o = \dots$

### Exercise 1: Use an analytical solution to formulate a 1D test

This exercise explores the exact solution (37). We shall formulate a diffusion problem in half of the domain for half of the Gaussian pulse. Then we shall investigate the impact of using an incorrect boundary condition, which we in general cases often are forced due if the solution needs to pass through finite boundaries undisturbed.

a) The solution (37) is seen to be symmetric at  $x = c$ , because  $\partial u / \partial x = 0$  always vanishes for  $x = c$ . Use this property to formulate a complete initial boundary value problem in 1D involving the diffusion equation  $u_t = \alpha u_{xx}$  on  $[0, L]$  with  $u_x(0, t) = 0$  and  $u(L, t)$  known.

b) Use the exact solution to set up a convergence rate test for an implementation of the problem. Investigate if a one-sided difference for  $u_x(0, t)$ , say  $u_0 = u_1$ , destroys the second-order accuracy in space.

c) Imagine that we want to solve the problem numerically on  $[0, L]$ , but we do not know the exact solution and cannot of that reason assign a correct Dirichlet condition at  $x = L$ . One idea is to simply set  $u(L, t) = 0$  since this will be an accurate approximation before the diffused pulse reaches  $x = L$  and even thereafter it might be a satisfactory condition. Let  $u_e$  be the exact solution and let  $u$  be the solution of  $u_t = \alpha u_{xx}$  with an initial Gaussian pulse and the boundary conditions  $u_x(0, t) = u(L, t) = 0$ . Derive a diffusion problem for the error  $e = u_e - u$ . Solve this problem numerically using an exact Dirichlet condition at  $x = L$ . Animate the evolution of the error and make a curve plot of the error measure

$$E(t) = \sqrt{\frac{\int_0^L e^2 dx}{\int_0^L u dx}}.$$

Is this a suitable error measure for the present problem?

d) Instead of using  $u(L, t) = 0$  as approximate boundary condition for letting the diffused Gaussian pulse out of our finite domain, one may try  $u_x(L, t) = 0$  since the solution for large  $t$  is quite flat. Argue that this condition gives a completely wrong asymptotic solution as  $t \rightarrow 0$ . To do this, integrate the diffusion equation from 0 to  $L$ , integrate  $u_{xx}$  by parts (or use Gauss' divergence theorem in 1D) to arrive at the important property

$$\frac{d}{dt} \int_0^L u(x, t) dx = 0,$$

implying that  $\int_0^L u dx$  must be constant in time, and therefore

$$\int_0^L u(x, t) dx = \int_0^L I(x) dx.$$

The integral of the initial pulse is 1.

e) Another idea for an artificial boundary condition at  $x = L$  is to use a cooling law

$$-\alpha u_x = q(u - u_S), \quad (52)$$

where  $q$  is an unknown heat transfer coefficient and  $u_S$  is the surrounding temperature in the medium outside of  $[0, L]$ . (Note that arguing that  $u_S$  is approximately  $u(L, t)$  gives the  $u_x = 0$  condition from the previous subexercise that is qualitatively wrong for large  $t$ .) Develop a diffusion problem for the error in the solution using (52) as boundary condition. Assume one can take  $u_S = 0$  "outside the domain" as  $u \rightarrow 0$  for  $x \rightarrow \infty$ . Find a function  $q = q(t)$  such that the exact solution obeys the condition (52). Test some constant values of  $q$  and animate how the corresponding error function behaves. Also compute  $E(t)$  curves as suggested in subexercise b). Filename: `diffu_symmetric_gaussian.py`.

### Exercise 2: Use an analytical solution to formulate a 2D test

Generalize (37) to multi dimensions by assuming that one-dimensional solutions can be multiplied to solve  $u_t = \alpha \nabla^2 u$ . Use this solution to formulate a 2D test case where the peak of the Gaussian is at the origin and where the domain is a rectangle in the first quadrant. Use symmetry boundary conditions  $\partial u / \partial n = 0$  wherever possible, and use exact Dirichlet conditions on the remaining boundaries. Filename: `diffu_symmetric_gaussian_2D.pdf`.

### Exercise 3: Examine stability of a diffusion model with a source term

Consider a diffusion equation with a linear  $u$  term:

$$u_t = \alpha u_{xx} + \beta u.$$

a) Derive in detail a Forward Euler scheme, a Backward Euler scheme, and a Crank-Nicolson for this type of diffusion model. Thereafter, formulate a  $\theta$ -rule to summarize the three schemes.

b) Assume a solution like (38) and find the relation between  $a$ ,  $k$ ,  $\alpha$ , and  $\beta$ .

c) Calculate the stability of the Forward Euler scheme. Design numerical experiments to confirm the results.

d) Repeat c) for the Backward Euler scheme.

e) Repeat c) for the Crank-Nicolson scheme.

f) How does the extra term  $bu$  impact the accuracy of the three schemes?

**Hint.** Compare the numerical and exact amplification factors, either in graphs or by Taylor series expansion (or both).

Filename: `diffu_stab_uterm.pdf`.

### 3 Exercises

#### Exercise 4: Stabilizing the Crank-Nicolson method by Rannacher time stepping

It is well known that the Crank-Nicolson method may give rise to non-physical oscillations in the solution of diffusion equations if the initial data exhibit jumps (see Section 2.6). Rannacher [1] suggested a stabilizing technique consisting of using the Backward Euler scheme for the first two time steps with step length  $\frac{1}{2}\Delta t$ . One can generalize this idea to taking  $2m$  time steps of size  $\frac{1}{2}\Delta t$  with the Backward Euler method and then continuing with the Crank-Nicolson method, which is of second-order in time. The idea is that the high frequencies of the initial solution are quickly damped out, and the Backward Euler scheme treats these high frequencies correctly. Thereafter, the high frequency content of the solution is gone and the Crank-Nicolson method will do well.

Test this idea for  $m = 1, 2, 3$  on a diffusion problem with a discontinuous initial condition. Measure the convergence rate using the solution (34) with the boundary conditions (35)-(36) for  $t$  values such that the conditions are in the vicinity of  $\pm 1$ . For example,  $t < 5a1.6 \cdot 10^{-2}$  makes the solution diffusion from a step to almost a straight line. The program `diffu_erf_sol.py` shows how to compute the analytical solution.

#### Project 5: Energy estimates for diffusion problems

This project concerns so-called *energy estimates* for diffusion problems that can be used for qualitative analytical insight and for verification of implementations.

a) We start with a 1D homogeneous diffusion equation with zero Dirichlet conditions:

$$u_t = \alpha u_x x, \quad x \in \Omega = (0, L), \quad t \in (0, T], \quad (53)$$

$$u(0, t) = u(L, t) = 0, \quad t \in (0, T], \quad (54)$$

$$u(x, 0) = I(x), \quad x \in [0, L]. \quad (55)$$

The energy estimate for this problem reads

$$\|u\|_{L^2} \leq \|I\|_{L^2}, \quad (56)$$

where the  $\|\cdot\|_{L^2}$  norm is defined by

$$\|g\|_{L^2} = \sqrt{\int_0^L g^2 dx}. \quad (57)$$

The quantity  $\|u\|_{L^2}$  or  $\frac{1}{2}\|u\|_{L^2}^2$  is known as the *energy* of the solution, although it is not the physical energy of the system. A mathematical tradition has introduced the notion *energy* in this context.

The estimate (56) says that the "size of  $u$ " never exceeds that of the initial condition, or more equivalently, that the area under the  $u$  curve decreases with time.

To show (56), multiply the PDE by  $u$  and integrate from 0 to  $L$ . Use that  $uu_t$  can be expressed as the time derivative of  $u^2$  and that  $u_x x u$  can be integrated by parts to form an integrand  $u_x^2$ . Show that the time derivative of  $\|u\|_{L^2}^2$  must be less than or equal to zero. Integrate this expression and derive (56).

b) Now we address a slightly different problem,

$$u_t = \alpha u_x x + f(x, t), \quad x \in \Omega = (0, L), \quad t \in (0, T], \quad (58)$$

$$u(0, t) = u(L, t) = 0, \quad t \in (0, T], \quad (59)$$

$$u(x, 0) = 0, \quad x \in [0, L]. \quad (60)$$

The associated energy estimate is

$$\|u\|_{L^2} \leq \|f\|_{L^2}. \quad (61)$$

(This result is more difficult to derive.)

Now consider the compound problem with an initial condition  $I(x)$  and a right-hand side  $f(x, t)$ :

$$u_t = \alpha u_x x + f(x, t), \quad x \in \Omega = (0, L), \quad t \in (0, T], \quad (62)$$

$$u(0, t) = u(L, t) = 0, \quad t \in (0, T], \quad (63)$$

$$u(x, 0) = I(x), \quad x \in [0, L]. \quad (64)$$

Show that if  $w_1$  fulfills (53)-(55) and  $w_2$  fulfills (58)-(60), then  $u = w_1 + w_2$  is the solution of (62)-(64). Using the triangle inequality for norms,

$$\|a + b\| \leq \|a\| + \|b\|,$$

show that the energy estimate for (62)-(64) becomes

$$\|u\|_{L^2} \leq \|I\|_{L^2} + \|f\|_{L^2}. \quad (65)$$

c) One application of (65) is to prove uniqueness of the solution. Suppose  $u_1$  and  $u_2$  both fulfill (62)-(64). Show that  $u = u_1 - u_2$  then fulfills (62)-(64) with  $f = 0$  and  $I = 0$ . Use (65) to deduce that the energy must be zero for all times and therefore that  $u_1 = u_2$ , which proves that the solution is unique.

d) Generalize (65) to a 2D/3D diffusion equation  $u_t = \nabla \cdot (\alpha \nabla u)$  for  $x \in \Omega$ .

**Hint.** Use integration by parts in multi dimensions:

$$\int_{\Omega} u \nabla \cdot (\alpha \nabla u) dx = - \int_{\Omega} \alpha \nabla u \cdot \nabla u dx + \int_{\partial\Omega} u \alpha \frac{\partial u}{\partial n},$$

where  $\frac{\partial u}{\partial n} = \mathbf{n} \cdot \nabla u$ ,  $\mathbf{n}$  being the outward unit normal to the boundary  $\partial\Omega$  of the domain  $\Omega$ .

e) Now we also consider the multi-dimensional PDE  $u_t = \nabla \cdot (\alpha \nabla u)$ . Integrate both sides over  $\Omega$  and use Gauss' divergence theorem,  $\int_{\Omega} \nabla \cdot \mathbf{q} \, dx = \int_{\partial\Omega} \mathbf{q} \cdot \mathbf{n} \, ds$  for a vector field  $\mathbf{q}$ . Show that if we have homogeneous Neumann conditions on the boundary,  $\partial u / \partial n = 0$ , area under the  $u$  surface remains constant in time and

$$\int_{\Omega} u \, dx = \int_{\Omega} I \, dx. \quad (66)$$

f) Establish a code in 1D, 2D, or 3D that can solve a diffusion equation with a source term  $f$ , initial condition  $I$ , and zero Dirichlet or Neumann conditions on the whole boundary.

We can use (65) and (66) as a partial verification of the code. Choose some functions  $f$  and  $I$  and check that (65) is obeyed at any time when zero Dirichlet conditions are used. Iterate over the same  $I$  functions and check that (66) is fulfilled when using zero Neumann conditions.

g) Make a list of some possible bugs in the code, such as indexing errors in arrays, failure to set the correct boundary conditions, evaluation of a term at a wrong time level, and similar. For each of the bugs, see if the verification tests from the previous subexercise pass or fail. This investigation shows how strong the energy estimates and the estimate (66) are for pointing out errors in the implementation.

Filename: `diffu_energy.pdf`.

## References

- [1] R. Rannacher. Finite element solution of diffusion problems with irregular data. *Numerische Mathematik*, 43:309–327, 1984.

## Index

amplification factor, 13

energy estimates (diffusion), 19

explicit discretization methods, 4

implicit discretization methods, 5

stationary solution, 3